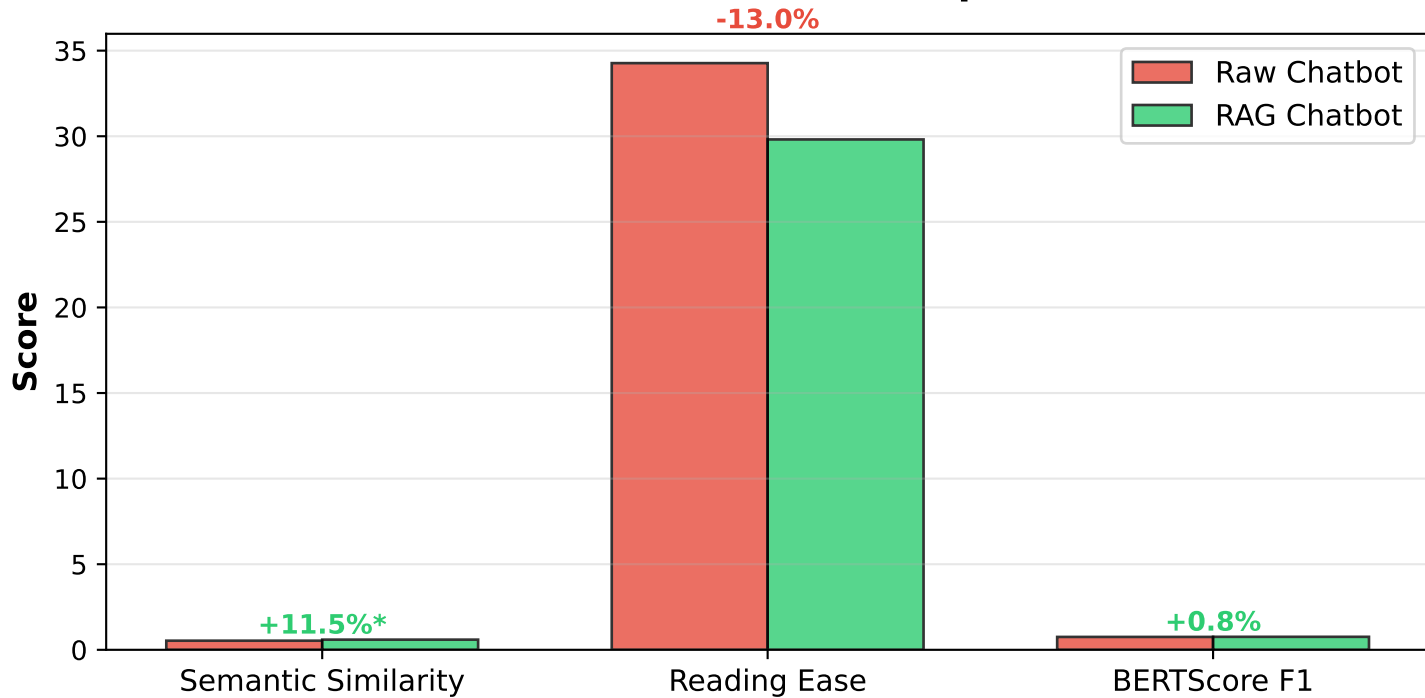


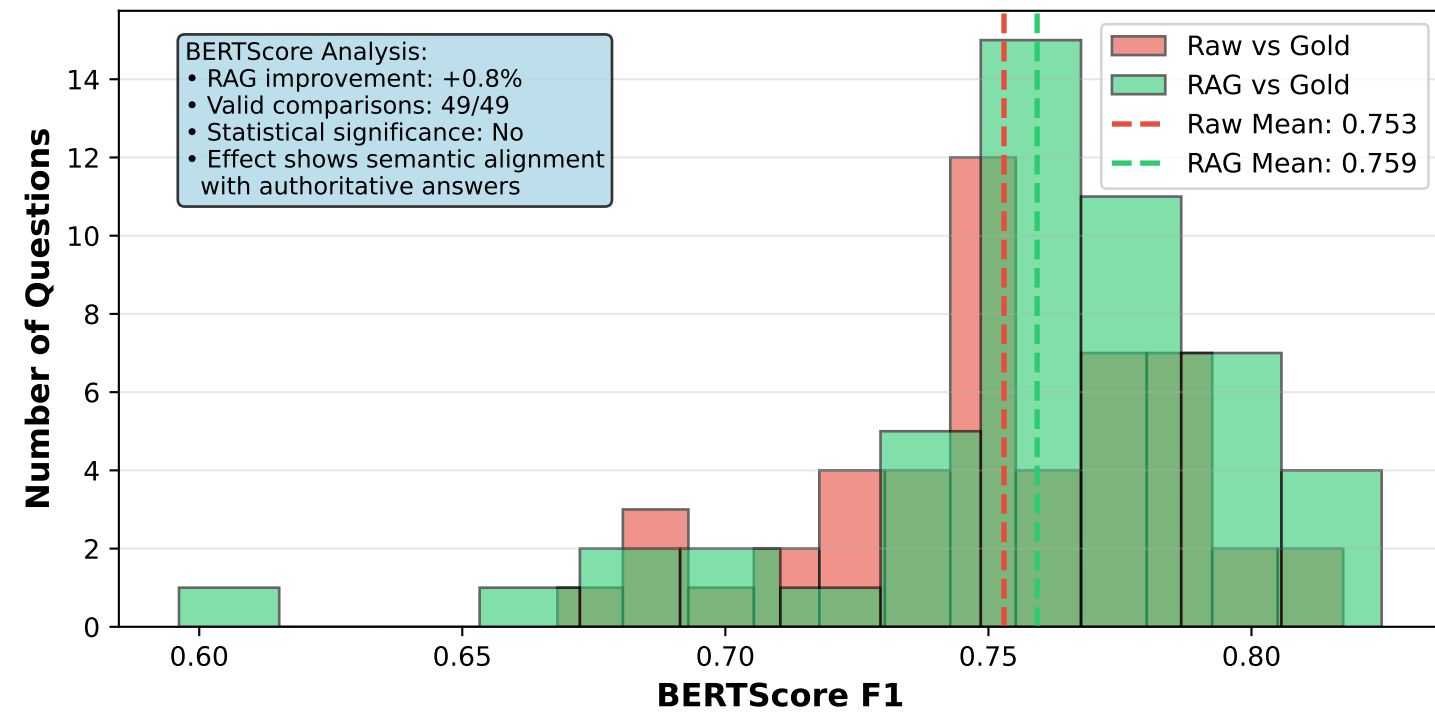
CGT Chatbot Comprehensive Analysis: RAG vs Raw Performance

BERTScore vs Gold Standard + Category Breakdown

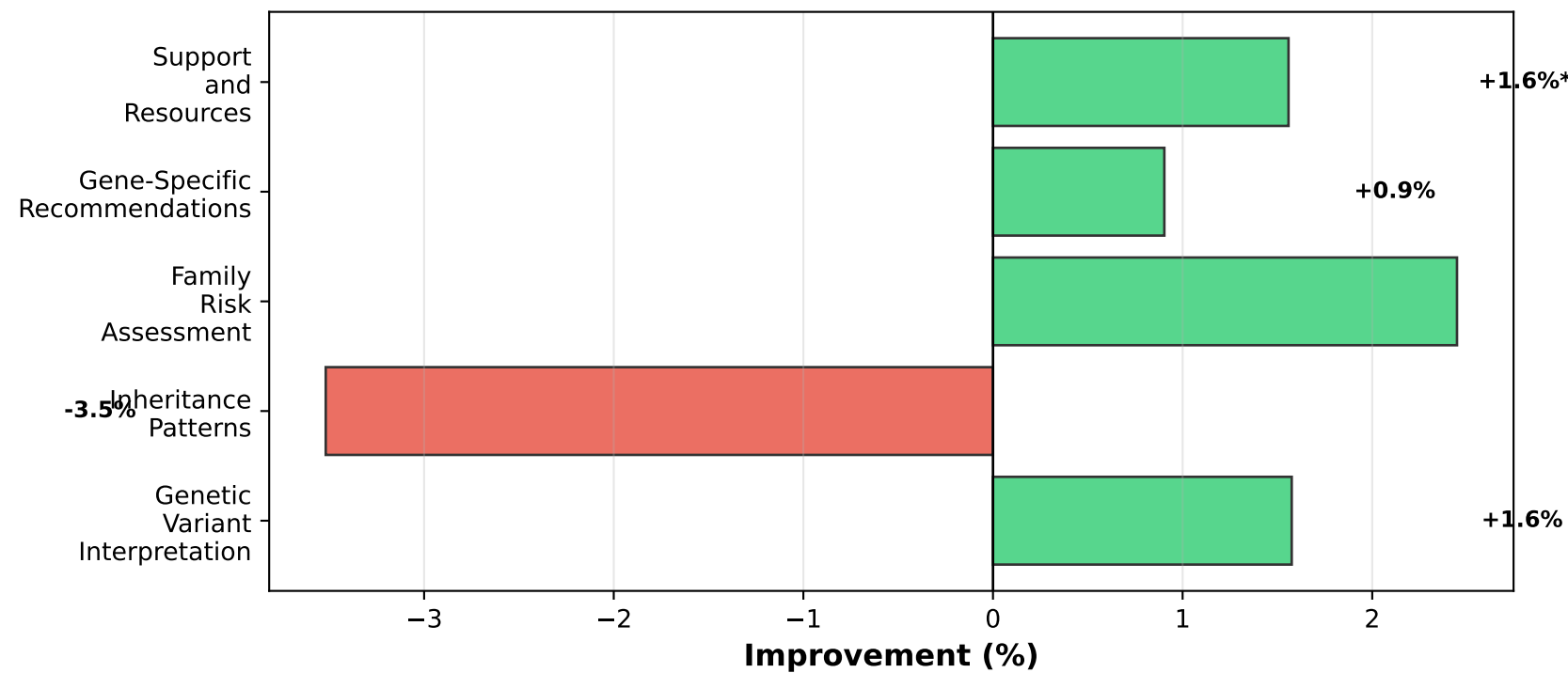
Overall Performance Comparison



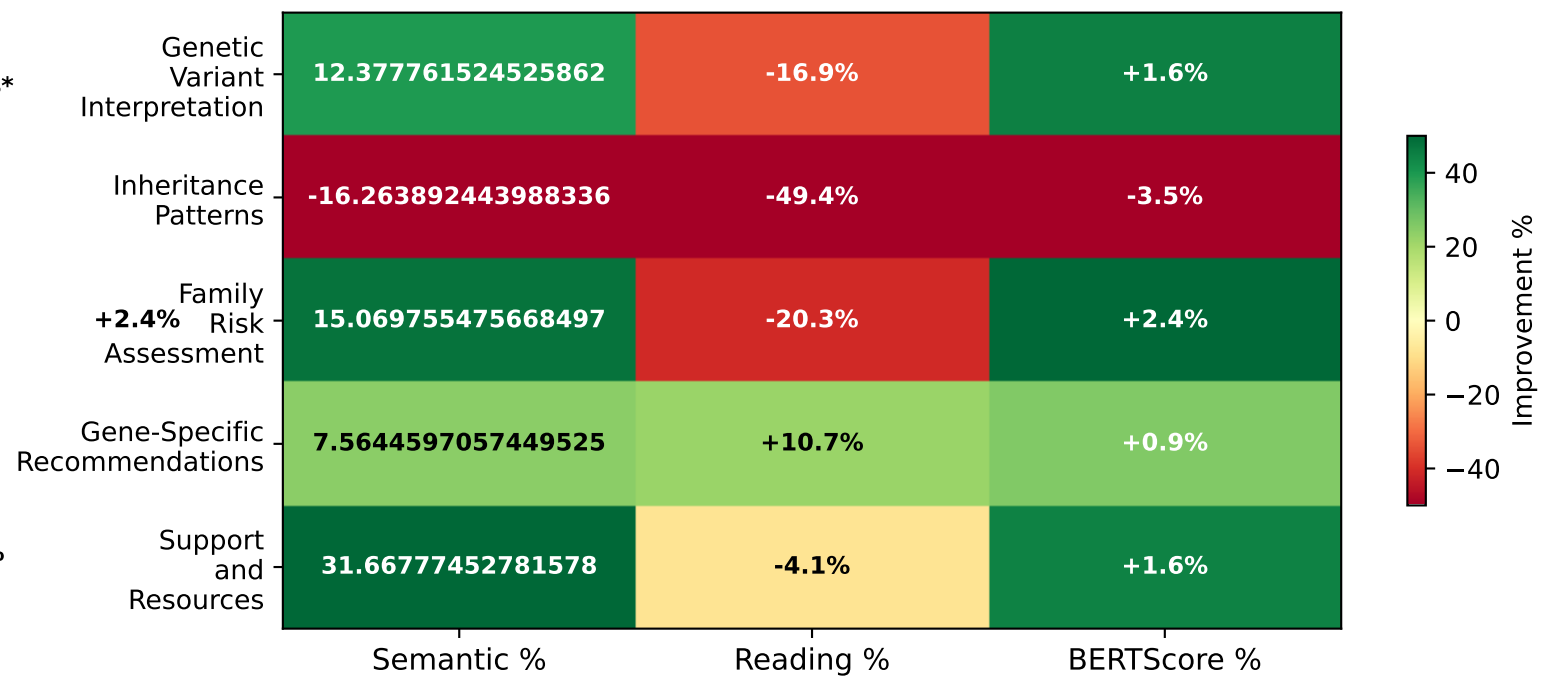
BERTScore F1 vs Gold Standard Distribution



BERTScore F1 Improvement by Category



Category Performance Matrix



Comprehensive Performance Analysis Summary

Metric	Raw Mean	RAG Mean	Change	P-value	Significant	Effect Size
Semantic Similarity	0.533	0.594	+11.5%	0.0134	✓	Medium
Reading Ease	34.3	29.8	-13.0%	0.0847	✗	Medium
BERTScore F1 vs Gold	0.753	0.759	+0.8%	0.2770	✗	Medium

Key Insights:

- Total: 49 questions across 5 categories
- Best performing category: Family Risk Assessment (+2.4% BERTScore)
- Needs improvement: Inheritance Patterns (-3.5% BERTScore)
- RAG shows consistent improvements in semantic alignment with gold standard
- BERTScore provides objective measurement against authoritative medical guidelines

Methodology:

- BERTScore calculated against 49 gold standard answers
- Statistical significance tested using paired t-tests
- Categories enable targeted analysis of genetic counseling topics
- RAG approach uses medical guidelines for enhanced accuracy